

CLAIMS

What is claimed is:

Sub
AI

1. A method of automatically creating a dictionary for clustering text documents comprising:
 - 5 determining a frequency of each word in each of said documents;
 - creating a Hashtable of most frequently occurring words in said documents;
 - determining a frequency of phrases in each of said documents that contain only words in said Hashtable;
 - adding most frequently occurring phrases to said Hashtable; and
 - 10 outputting said most frequently occurring words and said most frequently occurring phrases as said dictionary.
2. The method in claim 1, wherein said determining a frequency of each word comprises:
 - removing punctuation and case from said documents;
 - 15 removing stop words from said document;
 - replacing words in said documents with synonyms;
 - removing duplicate words from said documents;
 - adding remaining words to said Hashtable;
 - determining said frequency of each word remaining in said Hashtable; and
 - 20 removing words below a frequency level from said Hashtable.

3. The method in claim 2, further comprising inputting one or more of said stop words, said synonyms, and said frequency level.

4. The method in claim 1, wherein said determining a frequency of phrases comprises:

- 5 removing punctuation and case from said documents;
removing stop words from said document;
replacing words in said documents with synonyms;
adding said phrases in each of said documents that contain only words in
said Hashtable to said Hashtable;
10 determining said frequency of said phrases remaining in said Hashtable;
and
removing phrases below a frequency level from said Hashtable.

5. The method in claim 4, further comprising inputting one or more of said stop words, said synonyms, and said frequency level.

15 6. A method of automatically creating a dictionary for clustering text documents comprising:

- performing a first pass for each of said documents comprising:
determining a frequency of each word in each of said documents;
and

creating a Hashtable of most frequently occurring words in said documents;
performing a second pass for each of said documents comprising:
determining a frequency of phrases in each of said documents that
5 contain only words in said Hashtable; and
adding most frequently occurring phrases to said Hashtable; and
outputting said most frequently occurring words and said most frequently occurring phrases as said dictionary.

7. The method in claim 6, wherein said determining a frequency of each
10 word comprises:

removing punctuation and case from said documents;
removing stop words from said document;
replacing words in said documents with synonyms;
removing duplicate words from said documents;
15 adding remaining words to said Hashtable;
determining said frequency of each word remaining in said Hashtable; and
removing words below a frequency level from said Hashtable.

8. The method in claim 7, further comprising inputting one or more of said stop words, said synonyms, and said frequency level.

9. The method in claim 6, wherein said determining a frequency of phrases comprises:

removing punctuation and case from said documents;

removing stop words from said document;

5 replacing words in said documents with synonyms;

adding said phrases in each of said documents that contain only words in said Hashtable to said Hashtable;

determining said frequency of said phrases remaining in said Hashtable;

and

10 removing phrases below a frequency level from said Hashtable.

10. The method in claim 9, further comprising inputting one or more of said stop words, said synonyms, and said frequency level.

11. A program storage device readable by machine, tangibly embodying a

program of instructions executable by the machine to perform a method of

15 automatically creating a dictionary for clustering text documents, said method comprising:

determining a frequency of each word in each of said documents;

creating a Hashtable of most frequently occurring words in said documents;

15 determining a frequency of phrases in each of said documents that contain only words in said Hashtable;

20 adding most frequently occurring phrases to said Hashtable; and

outputting said most frequently occurring words and said most frequently occurring phrases as said dictionary.

12. A program storage device as in claim 11, wherein said determining a frequency of each word comprises:

- 5 removing punctuation and case from said documents;
- removing stop words from said document;
- replacing words in said documents with synonyms;
- removing duplicate words from said documents;
- adding remaining words to said Hashtable;
- 10 determining said frequency of each word remaining in said Hashtable; and
- removing words below a frequency level from said Hashtable.

13. A program storage device as in claim 12, further comprising inputting one or more of said stop words, said synonyms, and said frequency level.

14. A program storage device as in claim 11, wherein said determining a frequency of phrases comprises:

- 15 removing punctuation and case from said documents;
- removing stop words from said document;
- replacing words in said documents with synonyms;
- adding said phrases in each of said documents that contain only words in
- 20 said Hashtable to said Hashtable;

determining said frequency of said phrases remaining in said Hashtable;
and

removing phrases below a frequency level from said Hashtable.

5 15. A program storage device as in claim 14, further comprising inputting said stop words.

16. A program storage device as in claim 14, further comprising inputting said synonyms.

17. A program storage device as in claim 14, further comprising inputting said frequency level.